

(19) World Intellectual Property Organization  
International Bureau



(43) International Publication Date  
15 March 2001 (15.03.2001)

PCT

(10) International Publication Number  
**WO 01/19019 A1**

(51) International Patent Classification<sup>7</sup>: **H04L 12/00**

2071 (AU). BANH, Bui, Anh, Jonathan [AU/AU]; 34 Meryla Street, Burwood, NSW 2134 (AU).

(21) International Application Number: **PCT/AU00/01023**

(22) International Filing Date: 30 August 2000 (30.08.2000)

(74) Agent: **CONRICK, Patrick, Michael**; Intellectual Property Dept., Alcatel Australia Limited, 280 Botany Road, Alexandria, NSW 2015 (AU).

(25) Filing Language: **English**

(26) Publication Language: **English**

(81) Designated States (*national*): AU, CN, US.

(30) Priority Data:  
47400/99 6 September 1999 (06.09.1999) AU

(84) Designated States (*regional*): European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE).

(71) Applicant (*for all designated States except US*): **ALCATEL [FR/FR]**; 54, rue la Boétie, F-75009 Paris (FR).

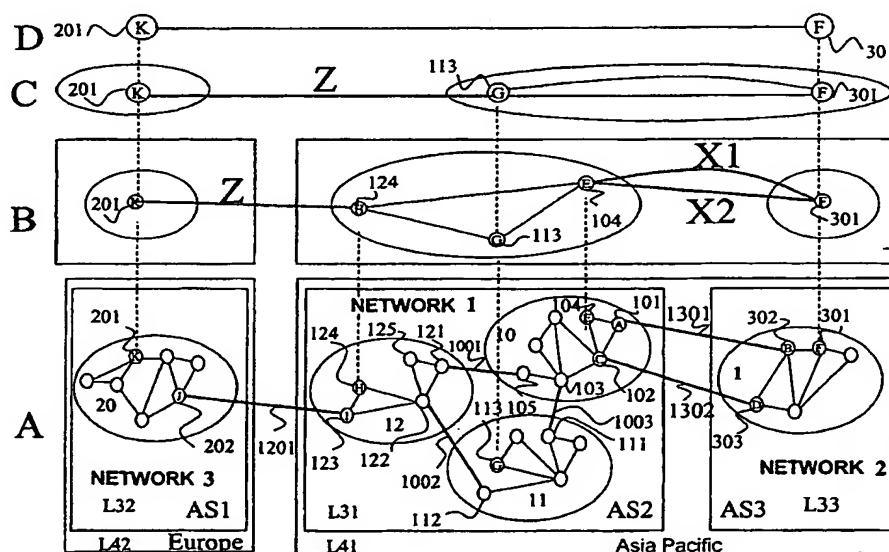
Published:  
— With international search report.

(72) Inventors; and

(75) Inventors/Applicants (*for US only*): **COX, Michael, S.** [AU/AU]; 101 7th Avenue, Jannali, NSW 2226 (AU). **VUCIC, Mickey** [AU/AU]; 25 Saiala Road, Killara, NSW

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(54) Title: **RECURSIVE TRAFFIC DISTRIBUTION IP/DATA NETWORK MODEL**



(57) Abstract: To avoid the need for all the nodes of a network to know the load status of all other nodes and links in the network, the nodes (101....303) are organized in node groups (10, 11, 12, 20, 30). Each group has a master node (104, 124, 113, 201, 301) which is incorporated in a logical higher order group. By recursively grouping the units of each preceding group, a logical hierarchical structure is formed in which the number of units at each level (A, B, C, D), decreases. Each node informs the other nodes of its load status, and each higher level unit also exchanges information with the other units. The recursive, hierarchical structure greatly reduces the amount of load status information exchanged across the network.

WO 01/19019 A1

RECURSIVE TRAFFIC DISTRIBUTION IP/DATA NETWORK MODELTechnical field

This invention relates to a method and arrangement for transferring network node status information between nodes.

5 Background Art

Least cost routing examines the network topography to determine the shortest path between the source of a message and its destination. This method does not take account of the load status of the individual links of the chosen path so one or more of the links may become overloaded, preventing or disrupting the delivery of the message to the destination.

10 An alternative proposal is to take into account the load status of the links when determining the chosen path. Thus the system may determine all the shortest paths and make the path selection on the basis of the path with the links carrying the least traffic. To implement this, it is necessary for all the nodes to exchange load status information. As the number of nodes and links in a network grow, the implementation of this technique requires the exchange of a large amount of load status traffic, as each node broadcasts the load status of its associated links to the other nodes.

15 In our co-pending application number 44470/99 (127045 SY), we disclose a network in which the message is spread over all the practical paths between the source and destination. An advantageous embodiment of that invention involves the elimination of paths having heavily loaded links. Again, the implementation of this technique requires the exchange of load status information between nodes, generating a large volume of traffic.

25 Disclosure of the Invention

This specification discloses a network arrangement for a plurality of nodes each node being connected to one or more other nodes by corresponding node links,

- 30 - the network being arranged into a recursive hierarchy of units having two or more levels,
- the nodes being the units of the first level of the hierarchy,

the units of higher levels of the hierarchy being formed by groupings of the units of the previous level,

wherein the units of a level exchange a corresponding load status information.

5 This invention can be used in conjunction with the maximal flow techniques described in 44470/99 to determine suitable paths with available capacity.

#### Brief Description of the Drawings

Figure 1 is a schematic representation of a network in which the nodes are arranged in a recursive hierarchy, in accordance with an embodiment of the

10 invention.

Figure 2 represents a load status monitor.

Figure 3 illustrates the message structure for exchanging information at different levels.

#### Best Mode of Carrying out the Invention

15 Figure 1 shows a network of nodes interconnected by links. According to the embodiment shown in Figure 1, the nodes are linked in a logical hierarchy.

The nodes exemplified by the circles, some of which are numbered 101 ... 202 ... 303 ..., are shown interconnected by node links, represented by the lines drawn between the nodes.

20 The nodes are formed into groups 10, 11, 12, 20, 30. The groups are interconnected by group links, for example 1001 between node 105 of group 10 and node 121 of group 12.

The group links are shown in the following Table 1.

TABLE 1 : GROUP LINKS

Node/Group		Node/Group		Link
105	10	121	12	1001
112	11	122	12	1002
103	10	111	11	1003
123	12	202	20	1201
101	10	302	30	1301
102	10	303	30	1302

Preferably the group links have a larger traffic capacity than node links. Group links may equate to regional trunks within a particular carrier's network, or to links between different carriers, different countries or different global regions; for example.

5 As shown in Figure 1, the levels of the hierarchy start with the nodes. The units of the second level are the groups of nodes. The units L31, L32, L33, of the third level are one or more groups, and the units of the fourth level, L41, L42 are formed by aggregating units of the third level. Thus L41 is the aggregation of L31 and L33, while L42 encompasses L32.

10 Preferably the nodes are grouped on the basis of communication path topography so that there are relatively few links in the shortest path between any two nodes in a group.

The groups themselves, 10, 11; 12, 20, 30 are linked by designated links 1001, 1002, 1003 1201, 1301, 1302 connecting designated nodes in the  
15 corresponding groups.

Within each group, a master node is assigned, 104, 113, 124, 201, 301. Because each node in a group knows the load status of all the links in that group, the master node has the information to enable it to compile available capacity reports adapted to meet the requirements of the higher levels.

20 The physical node interconnections are illustrated at A in Figure 1 and B, C and D illustrate the conceptual logical links for the second, third and fourth levels of hierarchy.

At level A, the nodes of each group communicate load status information to each of the other nodes of the corresponding group.

25 Thus the master nodes 104, 113, 124, of the network 1 shown at B are aggregated under the supervision of a single master node 113 which is assigned to the next higher level. Because, in the example shown, networks 2 and 3 have only one group, the same master 201, 301 is used all levels.

30 At level C, the master nodes 113, 201, 301 exchange information as to the overall load status of their associated networks 1, 2, and 3.

The nodes at level C are then 113, 201, 301. By then aggregating 113 and 301, one of these two can be designated to manage the level D information exchange for both network 1 and network 2. Thus at level D, 201 and 301 can exchange information on the available capacity for the regions covered by L42 and

5 L41:  
At level C, information on the capacity of networks 1, 2 and 3 is exchanged between the networks. At level B, information on the capacity of groups 10, 11, 12, 20, 30 is exchanged between the groups.

At level A, the nodes interchange information on traffic capacity at the node  
10 link level, within the groups.

Figure 2 shows an arrangement for monitoring the available capacity of the links connected to a node.

For each link connected to a node there is a buffer 51, 52, 53, e.g. in the form of FIFO.

15 The traffic level monitor 50 checks the level of the contents of the buffers to measure the available capacity on the basis of the speed of the link associated with the buffer. The result of the monitoring is then reported to the other nodes in the same group.

In a simplified measuring system, the monitor may report whether or not a  
20 link has spare capacity, e.g. by checking whether a buffer's content is above or below a predetermined threshold.

The load status information exchange is carried out on the following basis.

The nodes within a group each notify the other nodes within that group of the load status of the links connected to the notifying node.

25 At the group level, each group notifies the other groups of the load status of the links connected to the notifying group and a summary of the load status of internal paths within the group available for interconnecting the group links connected to the notifying group:

For example, Group 12 is connected to Group 20 via link 1201, to Group 10  
30 via link 1001, and to Group 11 via link 1002.

Preferably, the designated as a master node manages the interchange of information between the groups.

Table 2 shows the master nodes for each group.

TABLE 2

5

GROUP	MASTER NODE
10	104
11	113
12	124
20	201
30	301

As can be seen in Figure 1, the master nodes take part in the higher level exchanges but their number is progressively reduced by the recursive grouping.

Thus, in the embodiment shown in Figure 1, while there are 5 master nodes shown at level B, there are only 3 at level C and 2 at level D.

Preferably, the grouping is carried out on the basis of proximity in the sense of the number of links in the path. Of course this is not a strict rule at the node level because the nodes at either end of a group link are joined by a single link, while there may be more than 2 links between nodes within a group. Other factors which influence grouping are geographical proximity and network ownership, as well as the traffic flows.

For example, the nodes of network 2 may be geographically close to node of network 1, but network 1 may be owned by a different carrier from network 2.

At level D, the nodes 201 and 301 exchange information on the available capacity between network 3 and network 2 and the transit capacity of the respective networks. This information would, for example, be based on the load status of links 1201, 1301, 1302, and the capacity across network 1 between link 1201 and the links 1301, 1302. The information need only identify the maximum available capacity at the time, which varies in accordance with the load on the various network elements.

For the sake of clarity the information will be given the following names:

Level D = Regional;

Level C = Network;

Level B = Group;

5      Level A = Node.

Regional information may be, for example, the maximum available capacity between the "electrically" remotest groups. The term "electrically" refers to the number of links and may include cable, optical and radio links.

10      Network information may be, for example, the capacity between the various networks, including the trans-network capacity between the network links 1201, 1301, 1302.

Group information could be typified by the capacity between groups, including the trans-group capacity between the group links.

15      Node information is the information broadcast by a node to the other nodes within its group as the load status of the node and its associated links.

Group information can be deduced from node information. Each node in a group knows the load status of all the nodes in that group. Thus the master node 124 in group 12 knows the status of group-links 1001 from node 121, group link 1002 from node 122, and group/network link 1201 from node 123, as well as the  
20      status of all the internal nodes and links within group 12. Node 124 can therefore calculate the available capacity across the group 12 between any pair of the links 1201, 1001, 1002. Preferably the master node 124 would use the "all practical paths" algorithm of our Australian Patent application 44470/99 (Docket No. 127045 SY) to calculate the trans-group capacity. This group information is  
25      interchanged between the group master nodes 201, 124, 113, 104, 301 at level B.

The units of the level B group domain are again grouped together, in this embodiment, into 3 network groups. The network groups include two one member groups 201 and 301, and one three member group 124, 113, 104. The network master of each one member group is the member of the group, while 113 is  
30      designated as the master of the three member network group.

The three network masters from level B interchange network information at the level C network domain. The information relates to the network links connecting the respective networks, and the trans-network information relating to the capacity between the pairs of network links:

- 5        At the regional domain, level D, the network masters 201, 113, 301 have been formed into two groups, resulting in two regional masters 201, 301, which exchange information on the available capacity between the two regions.

10        The regional master nodes 201, 301, convey the regional link capacity information to the other regional nodes. In the present embodiment 301 conveys the information to 113. 201 is the only regional node in the other regional grouping.

The regional nodes 201, 113 and 301 are all network master nodes and they convey the inter-regional and inter-network capacity information to the network level nodes. In our embodiments, 113 conveys this information to the nodes 104, 124.

- 15        Each of the network level nodes 201, 124, 113, 104, 301 is a group master and relays the higher level information to each of the nodes in its group.

The grouping of the units at each level means that the information exchanged at each level becomes more generalised.

- 20        This means that a node has detailed capacity information about the other nodes in its group. Capacity information about other groups in its network, capacity information about the other networks in its region, and information about the inter-regional capacity.

- 25        In a preferred embodiment the group master handles the interchange of node link capacity information. Each node, instead of broadcasting its load status to all the other nodes in the group, sends the information only to the group master, which collates the information from each node and relays the information to the other nodes. The message from the group master preferably incorporates the higher level load status information, so that each node has an overall picture of the entire system.

- 30        Thus the group master may broadcast a message including the information shown in Figure 3. The first segment RL includes the load status at the regional link level D. A second portion of the payload includes a number of segments of



information on the inter-network load status NL. A third portion includes segments GL on the inter-group load status, and the fourth portion includes segments NL on the load status of the nodes within the group. Alternatively, this information can be flooded to other part of the network using other means, such as a broadcast or

## 5 multicast mechanism.

FIG. 10 is a block diagram of a network system 1000 according to one embodiment of the present invention. The network system 1000 includes a plurality of nodes 1010, 1020, 1030, 1040, 1050, 1060, 1070, 1080, 1090, 1100, 1110, 1120, 1130, 1140, 1150, 1160, 1170, 1180, 1190, 1200, 1210, 1220, 1230, 1240, 1250, 1260, 1270, 1280, 1290, 1300, 1310, 1320, 1330, 1340, 1350, 1360, 1370, 1380, 1390, 1400, 1410, 1420, 1430, 1440, 1450, 1460, 1470, 1480, 1490, 1500, 1510, 1520, 1530, 1540, 1550, 1560, 1570, 1580, 1590, 1600, 1610, 1620, 1630, 1640, 1650, 1660, 1670, 1680, 1690, 1700, 1710, 1720, 1730, 1740, 1750, 1760, 1770, 1780, 1790, 1800, 1810, 1820, 1830, 1840, 1850, 1860, 1870, 1880, 1890, 1900, 1910, 1920, 1930, 1940, 1950, 1960, 1970, 1980, 1990, 2000, 2010, 2020, 2030, 2040, 2050, 2060, 2070, 2080, 2090, 2100, 2110, 2120, 2130, 2140, 2150, 2160, 2170, 2180, 2190, 2200, 2210, 2220, 2230, 2240, 2250, 2260, 2270, 2280, 2290, 2300, 2310, 2320, 2330, 2340, 2350, 2360, 2370, 2380, 2390, 2400, 2410, 2420, 2430, 2440, 2450, 2460, 2470, 2480, 2490, 2500, 2510, 2520, 2530, 2540, 2550, 2560, 2570, 2580, 2590, 2600, 2610, 2620, 2630, 2640, 2650, 2660, 2670, 2680, 2690, 2700, 2710, 2720, 2730, 2740, 2750, 2760, 2770, 2780, 2790, 2800, 2810, 2820, 2830, 2840, 2850, 2860, 2870, 2880, 2890, 2900, 2910, 2920, 2930, 2940, 2950, 2960, 2970, 2980, 2990, 3000, 3010, 3020, 3030, 3040, 3050, 3060, 3070, 3080, 3090, 3100, 3110, 3120, 3130, 3140, 3150, 3160, 3170, 3180, 3190, 3200, 3210, 3220, 3230, 3240, 3250, 3260, 3270, 3280, 3290, 3300, 3310, 3320, 3330, 3340, 3350, 3360, 3370, 3380, 3390, 3400, 3410, 3420, 3430, 3440, 3450, 3460, 3470, 3480, 3490, 3500, 3510, 3520, 3530, 3540, 3550, 3560, 3570, 3580, 3590, 3600, 3610, 3620, 3630, 3640, 3650, 3660, 3670, 3680, 3690, 3700, 3710, 3720, 3730, 3740, 3750, 3760, 3770, 3780, 3790, 3800, 3810, 3820, 3830, 3840, 3850, 3860, 3870, 3880, 3890, 3900, 3910, 3920, 3930, 3940, 3950, 3960, 3970, 3980, 3990, 4000, 4010, 4020, 4030, 4040, 4050, 4060, 4070, 4080, 4090, 4100, 4110, 4120, 4130, 4140, 4150, 4160, 4170, 4180, 4190, 4200, 4210, 4220, 4230, 4240, 4250, 4260, 4270, 4280, 4290, 4300, 4310, 4320, 4330, 4340, 4350, 4360, 4370, 4380, 4390, 4400, 4410, 4420, 4430, 4440, 4450, 4460, 4470, 4480, 4490, 4500, 4510, 4520, 4530, 4540, 4550, 4560, 4570, 4580, 4590, 4600, 4610, 4620, 4630, 4640, 4650, 4660, 4670, 4680, 4690, 4700, 4710, 4720, 4730, 4740, 4750, 4760, 4770, 4780, 4790, 4800, 4810, 4820, 4830, 4840, 4850, 4860, 4870, 4880, 4890, 4900, 4910, 4920, 4930, 4940, 4950, 4960, 4970, 4980, 4990, 5000, 5010, 5020, 5030, 5040, 5050, 5060, 5070, 5080, 5090, 5100, 5110, 5120, 5130, 5140, 5150, 5160, 5170, 5180, 5190, 5200, 5210, 5220, 5230, 5240, 5250, 5260, 5270, 5280, 5290, 5300, 5310, 5320, 5330, 5340, 5350, 5360, 5370, 5380, 5390, 5400, 5410, 5420, 5430, 5440, 5450, 5460, 5470, 5480, 5490, 5500, 5510, 5520, 5530, 5540, 5550, 5560, 5570, 5580, 5590, 5600, 5610, 5620, 5630, 5640, 5650, 5660, 5670, 5680, 5690, 5700, 5710, 5720, 5730, 5740, 5750, 5760, 5770, 5780, 5790, 5800, 5810, 5820, 5830, 5840, 5850, 5860, 5870, 5880, 5890, 5900, 5910, 5920, 5930, 5940, 5950, 5960, 5970, 5980, 5990, 6000, 6010, 6020, 6030, 6040, 6050, 6060, 6070, 6080, 6090, 6100, 6110, 6120, 6130, 6140, 6150, 6160, 6170, 6180, 6190, 6200, 6210, 6220, 6230, 6240, 6250, 6260, 6270, 6280, 6290, 6300, 6310, 6320, 6330, 6340, 6350, 6360, 6370, 6380, 6390, 6400, 6410, 6420, 6430, 6440, 6450, 6460, 6470, 6480, 6490, 6500, 6510, 6520, 6530, 6540, 6550, 6560, 6570, 6580, 6590, 6600, 6610, 6620, 6630, 6640, 6650, 6660, 6670, 6680, 6690, 6700, 6710, 6720, 6730, 6740, 6750, 6760, 6770, 6780, 6790, 6800, 6810, 6820, 6830, 6840, 6850, 6860, 6870, 6880, 6890, 6900, 6910, 6920, 6930, 6940, 6950, 6960, 6970, 6980, 6990, 7000, 7010, 7020, 7030, 7040, 7050, 7060, 7070, 7080, 7090, 7100, 7110, 7120, 7130, 7140, 7150, 7160, 7170, 7180, 7190, 7200, 7210, 7220, 7230, 7240, 7250, 7260, 7270, 7280, 7290, 7300, 7310, 7320, 7330, 7340, 7350, 7360, 7370, 7380, 7390, 7400, 7410, 7420, 7430, 7440, 7450, 7460, 7470, 7480, 7490, 7500, 7510, 7520, 7530, 7540, 7550, 7560, 7570, 7580, 7590, 7600, 7610, 7620, 7630, 7640, 7650, 7660, 7670, 7680, 7690, 7700, 7710, 7720, 7730, 7740, 7750, 7760, 7770, 7780, 7790, 7800, 7810, 7820, 7830, 7840, 7850, 7860, 7870, 7880, 7890, 7900, 7910, 7920, 7930, 7940, 7950, 7960, 7970, 7980, 7990, 8000, 8010, 8020, 8030, 8040, 8050, 8060, 8070, 8080, 8090, 8100, 8110, 8120, 8130, 8140, 8150, 8160, 8170, 8180, 8190, 8200, 8210, 8220, 8230, 8240, 8250, 8260, 8270, 8280, 8290, 8300, 8310, 8320, 8330, 8340, 8350, 8360, 8370, 8380, 8390, 8400, 8410, 8420, 8430, 8440, 8450, 8460, 8470, 8480, 8490, 8500, 8510, 8520, 8530, 8540, 8550, 8560, 8570, 8580, 8590, 8600, 8610, 8620, 8630, 8640, 8650, 8660, 8670, 8680, 8690, 8700, 8710, 8720, 8730, 8740, 8750, 8760, 8770, 8780, 8790, 8800, 8810, 8820, 8830, 8840, 8850, 8860, 8870, 8880, 8890, 8900, 8910, 8920, 8930, 8940, 8950, 8960, 8970, 8980, 8990, 9000, 9010, 9020, 9030, 9040, 9050, 9060, 9070, 9080, 9090, 9100, 9110, 9120, 9130, 9140, 9150, 9160, 9170, 9180, 9190, 9200, 9210, 9220, 9230, 9240, 9250, 9260, 9270, 9280, 9290, 9300, 9310, 9320, 9330, 9340, 9350, 9360, 9370, 9380, 9390, 9400, 9410, 9420, 9430, 9440, 9450, 9460, 9470, 9480, 9490, 9500, 9510, 9520, 9530, 9540, 9550, 9560, 9570, 9580, 9590, 9600, 9610, 9620, 9630, 9640, 9650, 9660, 9670, 9680, 9690, 9700, 9710, 9720, 9730, 9740, 9750, 9760, 9770, 9780, 9790, 9800, 9810, 9820, 9830, 9840, 9850, 9860, 9870, 9880, 9890, 9900, 9910, 9920, 9930, 9940, 9950, 9960, 9970, 9980, 9990, 10000, 10010, 10020, 10030, 10040, 10050, 10060, 10070, 10080, 10090, 10100, 10110, 10120, 10130, 10140, 10150, 10160, 10170, 10180, 10190, 10200, 10210, 10220, 10230, 10240, 10250, 10260, 10270, 10280, 10290, 10300, 10310, 10320, 10330, 10340, 10350, 10360, 10370, 10380, 10390, 10400, 10410, 10420, 10430, 10440, 10450, 10460, 10470, 10480, 10490, 10500, 10510, 10520, 10530, 10540, 10550, 10560, 10570, 10580, 10590, 10600, 10610, 10620, 10630, 10640, 10650, 10660, 10670, 10680, 10690, 10700, 10710, 10720, 10730, 10740, 10750, 10760, 10770, 10780, 10790, 10800, 10810, 10820, 10830, 10840, 10850, 10860, 10870, 10880, 10890, 10900, 10910, 10920, 10930, 10940, 10950, 10960, 10970, 10980, 10990, 11000, 11010, 11020, 11030, 11040, 11050, 11060, 11070, 11080, 11090, 11100, 11110, 11120, 11130, 11140, 11150, 11160, 11170, 11180, 11190, 11200, 11210, 11220, 11230, 11240, 11250, 11260, 11270, 11280, 11290, 11300, 11310, 11320, 11330, 11340, 11350, 11360, 11370, 11380, 11390, 11400, 11410, 11420, 11430, 11440, 11450, 11460, 11470, 11480, 11490, 11500, 11510, 11520, 11530, 11540, 11550, 11560, 11570, 11580, 11590, 11600, 11610, 11620, 11630, 11640, 11650, 11660, 11670, 11680, 11690, 11700, 11710, 11720, 11730, 11740, 11750, 11760, 11770, 11780, 11790, 11800, 11810, 11820, 11830, 11840, 11850, 11860, 11870, 11880, 11890, 11900, 11910, 11920, 11930, 11940, 11950, 11960, 11970, 11980, 11990, 12000, 12010, 12020, 12030, 12040, 12050, 12060, 12070, 12080, 12090, 12100, 12110, 12120, 12130, 12140, 12150, 12160, 12170, 12180, 12190, 12200, 12210, 12220, 12230, 12240, 12250, 12260, 12270, 12280, 12290, 12300, 12310, 12320, 12330, 12340, 12350, 12360, 12370, 12380, 12390, 12400, 12410, 12420, 12430, 12440, 12450, 12460, 12470, 12480, 12490, 12500, 12510, 12520, 12530, 12540, 12550, 12560, 12570, 12580, 12590, 12600, 12610, 12620, 12630, 12640, 12650, 12660, 12670, 12680, 12690, 12700, 12710, 12720, 12730, 12740, 12750, 12760, 12770, 12780, 12790, 12800, 12810, 12820, 12830, 12840, 12850, 12860, 12870, 12880, 12890, 12900, 12910, 12920, 12930, 12940, 12950, 12960, 12970, 12980, 12990, 13000, 13010, 13020, 13030, 13040, 13050, 13060, 13070, 13080, 13090, 13100, 13110, 13120, 13130, 13140, 13150, 13160, 13170, 13180, 13190, 13200, 13210, 13220, 13230, 13240, 13250, 13260, 13270, 13280, 13290, 13300, 13310, 13320, 13330, 13340, 13350, 13360, 13370, 13380, 13390, 13400, 13410, 13420, 13430, 13440, 13450, 13460, 13470, 13480, 13490, 13500, 13510, 13520, 13530, 13540, 13550, 13560, 13570, 13580, 13590, 13600, 13610, 13620, 13630, 13640, 13650, 13660, 13670, 13680, 13690, 13700, 13710, 13720, 13730, 13740, 13750, 13760, 13770, 13780, 13790, 13800, 13810, 13820, 13830, 13840, 13850, 13860, 13870, 13880, 13890, 13900, 13910, 13920, 13930, 13940, 13950, 13960, 13970, 13980, 13990, 14000, 14010, 14020, 14030, 14040, 14050, 14060, 14070, 14080, 14090, 14100, 14110, 14120, 14130, 14140, 14150, 14160, 14170, 14180, 14190, 14200, 14210, 14220, 14230, 14240, 14250, 14260, 14270, 14280, 14290, 14300, 14310, 14320, 14330, 14340, 14350, 14360, 14370, 14380, 14390, 14400, 14410, 14420, 14430, 14440, 14450, 14460, 14470, 14480, 14490, 14500, 14510, 14520, 14530, 14540, 14550, 14560, 14570, 14580, 14590, 14600, 14610, 14620, 14630, 14640, 14650, 14660, 14670, 14680, 14690, 14700, 14710, 14720, 14730, 14740, 14750, 14760, 14770, 14780, 14790, 14800, 14810, 14820, 14830, 14840, 14850, 14860, 14870, 14880, 14890, 14900, 14910, 14920, 14930, 14940, 14950, 14960, 14970, 14980, 14990, 15000, 15010, 15020, 15030, 15040, 15050, 15060, 15070, 15080, 15090, 15100, 15110, 15120, 15130, 15140, 15150, 15160, 15170, 15180, 15190, 15200, 15210, 15220, 15230, 15240, 15250, 15260, 15270, 15280, 15290, 15300, 15310, 15320, 15330, 15340, 15350, 15360, 15370, 15380, 15390, 15400, 15410, 15420, 15430, 15440, 15450, 15460, 15470, 15480, 15490, 15500, 15510, 15520, 15530, 15540, 15550, 15560, 15570, 15580, 15590, 15600, 15610, 15620, 15630, 15640, 15650, 15660, 15670, 15680, 15690, 15700, 15710, 15720, 15730, 15740, 15750, 15760, 15770, 15780, 15790, 15800, 15810, 15820, 15830, 15840, 15850, 15860, 15870, 15880, 15890, 15900, 15910, 15920, 15930, 15940, 15950, 15960, 15970, 15980, 15990, 16000, 16010, 16020, 16030, 16040, 16050, 16060, 16070, 16080, 16090, 16100, 16110, 16120, 16130, 16140, 16150, 16160, 16170, 16180, 16190, 16200, 16210, 16220, 16230, 16240, 16250, 16260, 16270, 16280, 16290, 16300, 16310, 16320, 16330, 16340, 16350, 16360, 16370, 16380, 16390, 16400, 16410, 16420, 16430, 16440, 16450, 16460, 16470, 16480, 16490, 16500, 16510, 16520, 16530, 16540, 16550, 16560, 16570, 16580, 16590, 16600, 16610, 16620, 16630, 16640, 16650, 16660, 16670, 16680, 16690, 16700, 16710, 16720, 16730, 16740, 16750, 16760, 16770, 16780, 16790, 16800, 16810, 16820, 16830, 16840, 16850, 16860, 16870, 16880, 16890, 16900, 16910, 16920, 16930, 16940, 16950, 16960, 16970, 16980, 16990, 17000, 17010, 17020, 17030, 17040, 17050, 17060, 17070, 17080, 17090, 17100, 17110, 17120, 17130, 17140, 17150, 17160, 17170, 17180, 17190, 17200, 17210, 17220, 17230, 17240, 17250, 17260, 17270, 17280, 17290, 17300, 17310, 17320, 17330, 17340, 17350, 17360, 17370, 17380, 17390, 17400, 17410, 17420, 17430, 17440, 17450, 17460, 17470, 17480, 17490, 17500, 17510, 17520, 17530, 17540, 17550, 17560, 17570, 17580, 17590, 17600, 17610, 17620, 17630, 17640, 17650, 17660, 17670, 17680, 17690, 17700, 17710, 17720, 17730, 17740, 17750, 17760, 17770, 17780, 17790, 17800, 17810, 17820, 17830, 17840, 17850, 17860, 17870, 17880, 17890, 17900, 17910, 17920, 17930, 17940, 17950, 17960, 17970, 17980, 17990, 18000, 18010, 18020, 18030, 18040, 18050, 18060, 18070, 18080, 18090, 18100, 18110, 18120, 18130, 18140, 18150, 18160, 18170, 18180, 18190, 18200, 18210, 18220, 18230, 18240, 18250, 18260, 18270, 18280, 18290, 18300, 18310, 18320, 18330, 18340, 18350, 18360, 18370, 18380, 18390, 18400, 18410, 18420, 18430, 18440, 18450, 18460, 18470, 18480, 18490, 18500, 18510, 18520, 18530, 18540, 18550, 18560, 18570, 18580, 18590, 18600, 18610, 18620, 18630, 18640, 18650, 18660, 18670, 18680, 18690, 18700, 18710, 18720, 18730, 18740, 18750, 18760, 18770, 18780, 18790, 18800, 18810, 18820, 18830, 18840, 18850, 18860, 18870, 18880, 18890, 18900, 18910, 18920, 18930, 18940, 18950, 18960, 18970, 18980, 18990, 19000, 19010, 19020, 19030, 19040, 19050, 19060, 19070, 19080, 19090, 19100, 19110, 19120, 19130, 19140, 19150, 19160, 19170, 19180, 19190, 19200, 19210, 19220, 19230, 19240, 19250, 19260, 19270, 19280, 19290, 19300, 19310, 19320, 19330, 19340, 19350, 19360, 19370, 19380, 19390, 19400, 19410, 19420, 19430, 19440, 19450, 19460, 19470, 19480, 19490, 19500, 19510, 1

The claims defining the invention are as follows:

1. A network arrangement for a plurality of nodes, each node being connected to one or more other nodes by corresponding node links,
  - the network being arranged into a recursive hierarchy of units having
- 5 two or more levels,
  - the nodes being the units of the first level of the hierarchy,
  - the units of higher levels of the hierarchy being formed by groupings of the units of the previous level,
  - wherein the units of a level exchange a corresponding load status
- 10 information.
2. An arrangement as claimed in claim 1 wherein, within, each group of units, a master entity is designated, the master entity conveying inter-unit load status information relating to the units of that level to the next higher level.
3. An arrangement as claimed claim 1 or claim 2 wherein, in the first level, a
- 15 selected node in each group is designated as the master node for the corresponding group,
  - the master node managing the transfer of node load status information within its corresponding group.
4. An arrangement as claimed in claim 1 or claim 2, or claim 3 wherein the
- 20 load status information includes information on the available traffic capacity between the ports of each unit.
5. An arrangement as claimed in any one of claims 1 to 4 wherein each node includes node load status monitoring means to monitor the load status of the links connected to the node.
- 25 6. An arrangement as claimed in any one of claims 1 to 5 wherein at least one node of each second level group is connected to a node of at least one other second level group via a corresponding group link whereby group load status information can be interchange.
7. An arrangement as claimed in claim 6 wherein the units of the third level are
- 30 formed by mutually interconnected second level units.

8. A network arrangement for interchanging load status information substantially as herein described with reference to the accompanying drawings.

9. A network arrangement as claimed in any one of claims 1 to 8 implementing the maximal flow techniques of 44470/99.

5

**This Page Blank (uspto)**

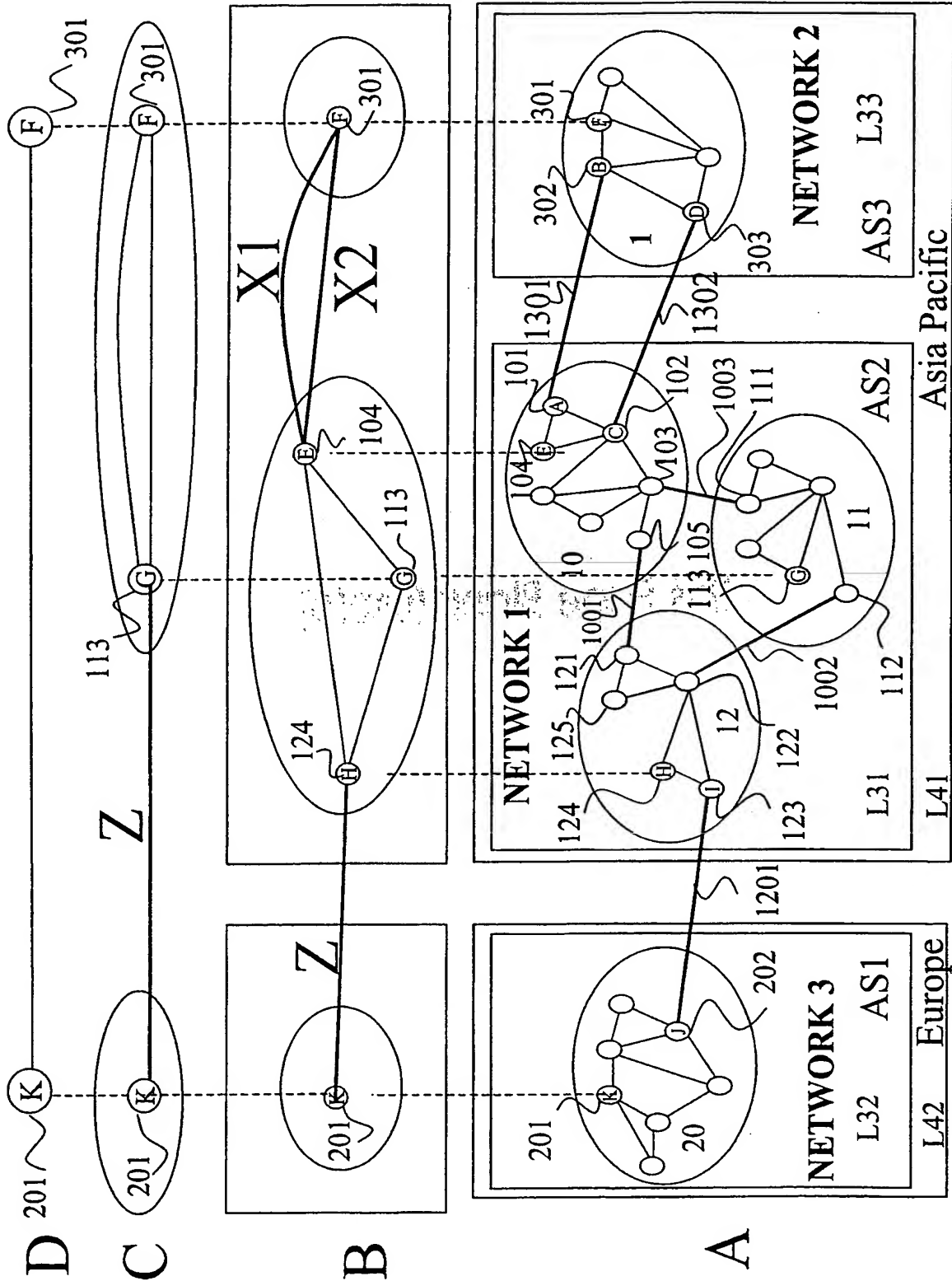


FIGURE 1

**This Page Blank (uspto)**

2/3

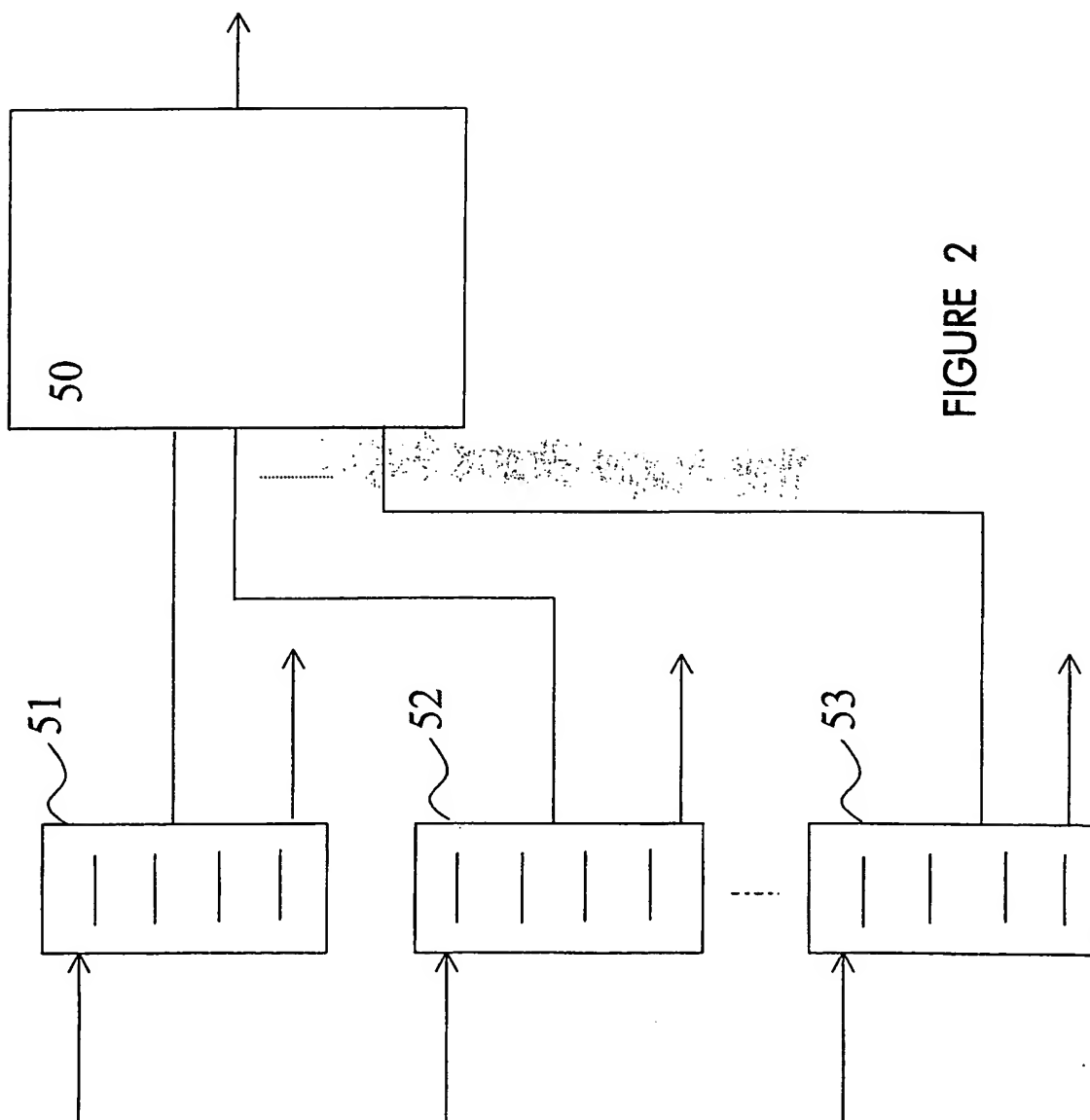


FIGURE 2

**This Page Blank (uspto)**



RL	NL	NL	NL	GL	GL	GL	GL	GL	GL	NL	NL	NL	NL	NL	NL
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

FIGURE 3

**This Page Blank (uspto)**

# INTERNATIONAL SEARCH REPORT

International application No.  
PCT/AU 00/01023

<b>A. CLASSIFICATION OF SUBJECT MATTER</b>		
Int Cl <sup>7</sup> : H04L 12/00		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b>		
Minimum documentation searched (classification system followed by classification symbols) WHOLE IPC		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched AU: IPC as above		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) WPAT: (NETWORK + or LAN or Wan) and (HIERARCH + or TREE or LAYER) and (DISTRIBUTION? or LOAD or TRAFFIC or FLOW or CONGESTION) and (STATUS or CONDITION or STATE) INSPEC: DITTO		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	US 5483533A (KUBA) 9 January 1996 Whole document	1-9
A	US 5537468A (HARTMANN) 16 July 1996 Whole document	1-9
A	US 5872773A (KATZELA et al.) 16 February 1999 Whole document	1-9
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C <span style="margin-left: 100px;"><input checked="" type="checkbox"/> See patent family annex</span>		
<p>* Special categories of cited documents:</p> <p>"A" Document defining the general state of the art which is not considered to be of particular relevance</p> <p>"E" earlier application or patent but published on or after the international filing date</p> <p>"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)</p> <p>"O" document referring to an oral disclosure, use, exhibition or other means</p> <p>"P" document published prior to the international filing date but later than the priority date claimed</p>	<p>"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention</p> <p>"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone</p> <p>"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art</p> <p>"&amp;" document member of the same patent family</p>	
Date of the actual completion of the international search 29 September 2000		Date of mailing of the international search report 12 OCT 2000
Name and mailing address of the ISA/AU AUSTRALIAN PATENT OFFICE PO BOX 200 WODEN ACT 2606 AUSTRALIA E-mail address: pct@ipaaustralia.gov.au Facsimile No.: (02) 6285 3929		Authorized officer  JUZER KHANBHAI Telephone No.: (02) 6283 2176

# INTERNATIONAL SEARCH REPORT

international application No.

PCT/AU 00/01023

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	DE 19742582 C1 (SIEMENS AG) 29 April 1999. Whole document	1-9
A	US 5905871A (BUSKENS et al.) 18 May 1999 Whole document	1-9

# INTERNATIONAL SEARCH REPORT

## Information on patent family members

International application No.  
PCT/AU 00/01023

This Annex lists the known "A" publication level patent family members relating to the patent documents cited in the above-mentioned international search report. The Australian Patent Office is in no way liable for these particulars which are merely given for the purpose of information.

Patent Document Cited in Search Report			Patent Family Member	
US	5483533	JP	7087112	
US	5537468	EP	608279	WO 9308666
US	5872773	NONE		
DE	19742582	NONE		
US	5905871	NONE		

**This Page Blank (uspto)**